

Lecture Notes in Computer Science

Edited by G. Goos and J. Hartmanis

205

Paul Klint

A Study in
String Processing Languages



Springer-Verlag
Berlin Heidelberg New York Tokyo

Editorial Board

D. Barstow W. Brauer P. Brinch Hansen D. Gries D. Luckham
C. Moler A. Pnueli G. Seegmüller J. Stoer N. Wirth

Author

Paul Klint
Centre for Mathematics and Computer Science
Kruislaan 413, 1098 SJ Amsterdam, The Netherlands

CR Subject Classification (1985): D.3.1, D.3.3, F.3.3, I.7.0

ISBN 3-540-16041-8 Springer-Verlag Berlin Heidelberg New York Tokyo
ISBN 0-387-16041-8 Springer-Verlag New York Heidelberg Berlin Tokyo

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically those of translation, reprinting, re-use of illustrations, broadcasting, reproduction by photocopying machine or similar means, and storage in data banks. Under § 54 of the German Copyright Law where copies are made for other than private use, a fee is payable to "Verwertungsgesellschaft Wort", Munich.

© by Springer-Verlag Berlin Heidelberg 1985
Printed in Germany

Printing and binding: Beltz Offsetdruck, Hemsbach/Bergstr.
2145/3140-543210

CONTENTS

CONTENTS *iii*

PREFACE *vii*

PART I: String Processing Languages

1. INTRODUCTION 3

- 1.1. Subject of this monograph 3
- 1.2. Basic operations on strings 5
- 1.3. Why are string processing languages special? 8
 - 1.3.1. Bookkeeping 8
 - 1.3.2. Recognition strategy 8
 - 1.3.3. Failure handling 9
 - 1.3.4. Existing languages and string processing 10
- 1.4. Problems in string processing languages 11
 - 1.4.1. A short introduction to SNOBOL4 11
 - 1.4.2. Compound patterns 12
 - 1.4.3. Side-effects during pattern matching 13
 - 1.4.4. Problems with the SNOBOL4 approach 14
- 1.5. A checklist for string processing languages 15
 - 1.5.1. Treatment of the subject 15
 - 1.5.2. Recognition strategy 15
- 1.6. References for Chapter 1 16

2. DESIGN CONSIDERATIONS FOR STRING PROCESSING LANGUAGES 17

- 2.1. Introduction 17
- 2.2. Some representative pattern matching functions and operators 18
- 2.3. Description methods for pattern matching 18
 - 2.3.1. Patterns defined by sets of strings 19
 - 2.3.2. Patterns defined by algebraic transformations 19
 - 2.3.3. Patterns defined by recursive coroutines 20
 - 2.3.4. Patterns defined by operational semantics 20

2.4. A comparison of two backtracking models	21
2.4.1. Common definitions for the two models	21
2.4.2. The immediate/conditional model	23
2.4.2.1. Overview	23
2.4.2.2. Formal description	25
2.4.3. The recovery model	31
2.4.3.1. Overview	31
2.4.3.2. Formal description	33
2.5. Unification of pattern and expression language	35
2.6. References for Chapter 2	36
3. AN OVERVIEW OF THE SUMMER PROGRAMMING LANGUAGE	38
3.1. Introduction	38
3.2. Success-directed evaluation and control structures	38
3.3. Recovery of side-effects	40
3.4. Procedures, operators and classes	42
3.5. A pattern matching extension	44
3.5.1. String Pattern Matching	44
3.5.2. Generalized pattern matching	47
3.6. Related work	48
3.7. References for Chapter 3	48
4. FORMAL LANGUAGE DEFINITIONS CAN BE MADE PRACTICAL	49
4.1. The problem	49
4.2. The method	50
4.2.1. Introduction	50
4.2.2. SUMMER as a metalanguage	50
4.2.3. Semantic domains	52
4.2.4. Evaluation process	53
4.2.5. Some examples	56
4.2.5.1. If expressions	56
4.2.5.2. Variable declarations	57
4.2.5.3. Blocks	58
4.3. Assessment	59
4.4. References for Chapter 4	61
5. ASSESSMENT	62
5.1. Looking backward	62
5.1.1. SUMMER as a language	62
5.1.2. The SUMMER implementation	63
5.1.3. Use of a formal definition	63
5.2. Looking forward	64
5.3. References for Chapter 5	65

PART II: Definition of SUMMER**6. PRELIMINARIES TO THE DEFINITION OF SUMMER 69**

- 6.1. Syntactic considerations 69
- 6.2. Lexical considerations 70
- 6.3. Semantic considerations 71
 - 6.3.1. Description method 71
 - 6.3.2. SUMMER as a metalanguage 72
 - 6.3.3. Semantic domains 74
 - 6.3.4. Evaluation process 78
- 6.4. Features not specified in the definition 81
- 6.5. References for chapter 6 82

7. A SEMI-FORMAL DEFINITION OF THE SUMMER KERNEL 83

- 7.1. Declarations 83
 - 7.1.1. Summer program 83
 - 7.1.2. Variable declarations 85
 - 7.1.3. Constant declarations 86
 - 7.1.4. Procedure and operator declarations 87
 - 7.1.5. Class declarations 88
 - 7.1.6. Operator symbol declarations 92
- 7.2. Expressions 93
 - 7.2.1. Constants 94
 - 7.2.2. Identifiers and procedure calls 94
 - 7.2.3. Return expressions 99
 - 7.2.4. If expressions 100
 - 7.2.5. Case expressions 103
 - 7.2.6. While expressions 104
 - 7.2.7. For expressions 105
 - 7.2.8. Try expressions 107
 - 7.2.9. Scan expressions 109
 - 7.2.10. Assert expressions 110
 - 7.2.11. Parenthesized expressions and blocks 111
 - 7.2.12. Array expressions 112
 - 7.2.13. Table expressions 115
 - 7.2.14. Field selection 117
 - 7.2.15. Subscription 120
 - 7.2.16. Monadic expressions 122
 - 7.2.17. Dyadic expressions 123
 - 7.2.18. Constant expressions 126
- 7.3. Miscellaneous functions used in the formal definition 128
 - 7.3.1. The function *dereference* 128
 - 7.3.2. The function *equal* 128
 - 7.3.3. The functions *substring* and *string_equal* 129

8. THE SUMMER LIBRARY 130

- 8.1. Introduction *130*
- 8.2. Class integer *130*
- 8.3. Class real *132*
- 8.4. Class string *134*
- 8.5. Class array *138*
- 8.6. Class interval *141*
- 8.7. Class table *142*
- 8.8. Class scan_string *143*
- 8.9. Class file *147*
- 8.10. Class bits *149*
- 8.11. Miscellaneous procedures *150*

9. SOME ANNOTATED SUMMER PROGRAMS 151

- 9.1. Introduction *151*
 - 9.2.1. Word tuples *151*
 - 9.2.2. Flexible arrays *154*

10. SUMMARY OF SUMMER SYNTAX 158**INDEX FOR PART II 161**

PREFACE

Written text is an essential element in our culture and various technical means have been invented to aid in its production. Paper and pencil, the typewriter and the typesetter are examples of such inventions.

Continuing this same line of development, computers are nowadays being used to alleviate the writing task. Computerized text processing systems (ranging from word processors for writing and editing simple texts to fully automated newspaper and book printing systems) are rapidly penetrating into all areas of human activity where written text is the primary means of communication.

Historically, the impetus behind the development of computers has always been primarily numerical in nature. This is reflected in the design of most computers and programming languages. However, the increasing use of computers for text processing and other non-numeric tasks makes the purely arithmetic design obsolete.

This monograph concentrates on the programming language aspects of computerized text handling and, to be more precise, on the design and implementation of **string processing languages**. The term 'string processing' refers to the process of inspecting, modifying and transforming texts, i.e. sequences of symbols. It comprises such seemingly disparate activities as text editing, transforming a text with embedded formatting directives into a final layout, and compiling a source program into a string of machine instructions.

An account is given of attempts to solve some of the problems in string processing languages. First of all, an exercise in designing an application oriented programming language is described. This has resulted in the language SUMMER. Next, an exercise in the formal definition of the semantics of programming languages is described. The formal definition and implementation of SUMMER together constitute the final result of the project.

This monograph consists of two parts. Part I is devoted to string processing languages in general and consists of chapters 1 through 5. Part II is devoted to the definition of SUMMER and consists of chapters 6 through 10. The contents of the monograph are now briefly summarized.

Chapter 1 is introductory and gives the necessary motivation and background for the study of string processing languages.

Chapter 2 is devoted to general design considerations for string processing languages and compares the semantics of various pattern matching models. Attention is paid to different forms of side-effects during a pattern match. This is done by giving an operational, formal definition of the semantics of the various models. As a result of this, a new pattern matching model based on side-effect recovery is developed.

Chapter 3 gives an overview of the language SUMMER. SUMMER may be characterized as a **small** language, i.e. it consists of a relatively small set of primitive operations together with a modest extension mechanism.

Chapter 4 concentrates on the problem of finding a method for formal language definition that is suitable for the designers as well as the implementors and users of a language. An improved method for the operational definition of programming language semantics is developed and the result of applying this method to SUMMER is illustrated.

Chapter 5 concludes the first part of this monograph with an evaluation of the research described in it and suggestions for further research.

Part II is devoted to the definition of the SUMMER programming language. It provides both a formal and informal language definition and tutorial examples.

In Chapter 6 the techniques and notational conventions that are used in the definition are introduced. Much attention is paid to the method used for the formal definition of the semantics of SUMMER.

Chapter 7 contains a semi-formal definition of the SUMMER kernel. This is a small subset of the language on which a semantic definition of the whole language can be based. The description of each language feature consists of its syntax, an informal as well as a formal definition of its semantics, and examples.

In Chapter 8 the kernel is extended with useful data types and associated operations, such as reals, arrays, tables, files, bit strings, etc. Some complete, annotated SUMMER programs are presented in Chapter 9. Finally, a summary of the syntax is given in Chapter 10.

Readers who are only interested in getting a general impression of the language SUMMER may confine themselves to Chapter 3 and the annotated examples in Chapter 9. Readers who are not interested in the formal definition of the language may skip Chapter 6 (except Sections 6.1 and 6.2), and all subsections of Chapter 7 entitled 'Semantics'.

This monograph is a revision of the author's dissertation as presented at the Technical University Eindhoven, The Netherlands. It was supervised by Prof. dr. F.E.J. Kruseman Aretz (Technical University Eindhoven, The Netherlands) and Prof. H. Whitfield, b.s., d.i.c. (University of Newcastle upon Tyne, Great Britain). The research was conducted while the author was employed at the Mathematical Centre, Amsterdam, The Netherlands.

Several people contributed to this effort. Design and implementation of SUMMER (as well as of its predecessor SPRING) were done in close cooperation with Marleen Sint and contributions to the design of SUMMER were made by Jan Heering. Their enthusiasm, patience and friendship were essential to the success of this project.

Jan Heering, Marleen Sint and Arthur Veen have read drafts of this monograph. They pointed out various errors and made numerous suggestions for improving the style and presentation of it. I am grateful for their support and criticism. Comments made by Leo Geurts, R.J. Lunbeck, Lambert Meertens and W.L. van der Poel are gratefully acknowledged.